

Painting by evolutionary history: Inference of local ancestry in admixed populations



Ali J. Berens (ali.berens@biology.gatech.edu), Joseph Lachance
School of Biology, Georgia Institute of Technology



Background

Local ancestry inference is important for:

- Understanding history of populations
- Personalized medicine

STRUCTURE-like approaches do not paint local ancestry

Most existing methods require reference sequences

Objective: Develop a reference-free, hierarchical local ancestry inference method

Test Data

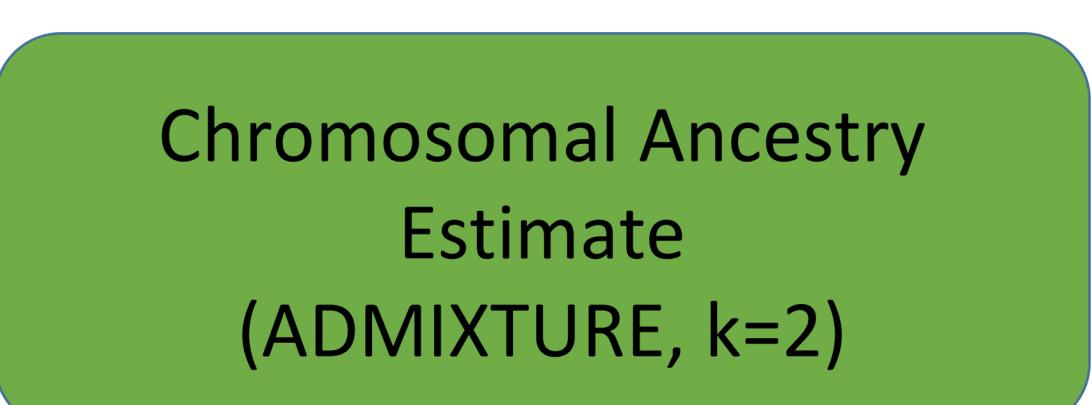
1000 Genomes Project: ASW, CEU, YRI populations

Phased chromosomes

SNP filtering: Biallelic and MAF > 0.05

Workflow

For full dataset:



Q-file	
A1	A2
Sample 1 (C1) 0.763	0.247
Sample 2 (C2) 0.155	0.845
⋮	⋮

Select Informative SNPs

Create haplotypes of l consecutive SNPs

Create emission matrix for each window of length l

Phased Chromosomes		Top frequency differences
C1	... Cj	A1 A2
SNP 1 [0 ... 1]	...	SNP 1 [0.10 0.15]
SNP i [1 ... 0]	...	SNP 2 [0.33 0.24]
⋮	⋮	⋮

Haplotypes ($l = 5$)

C1	... Cj
SNPs 1–5 [00101 ... 01100]	⋮
SNPs 10–15 [00101 ... 01100]	⋮
⋮	⋮

SNPs 1–5	haplotype sums	Emission matrix
C1 [00101]	weighted by Q-file	00000 ... 11111
Cj [01100]	A1 [0.674 ... 0.002]	A2 [0.622 ... 0.001]

For each sample:

Create transition matrix

$$\begin{array}{cc} \text{Previous State} & \text{Next State} \\ \begin{matrix} A1 \\ A2 \end{matrix} & \begin{matrix} A1 & A2 \\ rp_2 & 1 - rp_2 \\ rp_1 & 1 - rp_1 \end{matrix} \end{array}$$

r: recombination rate
 p_1, p_2 : from Q-file

Forward HMM

Backward HMM

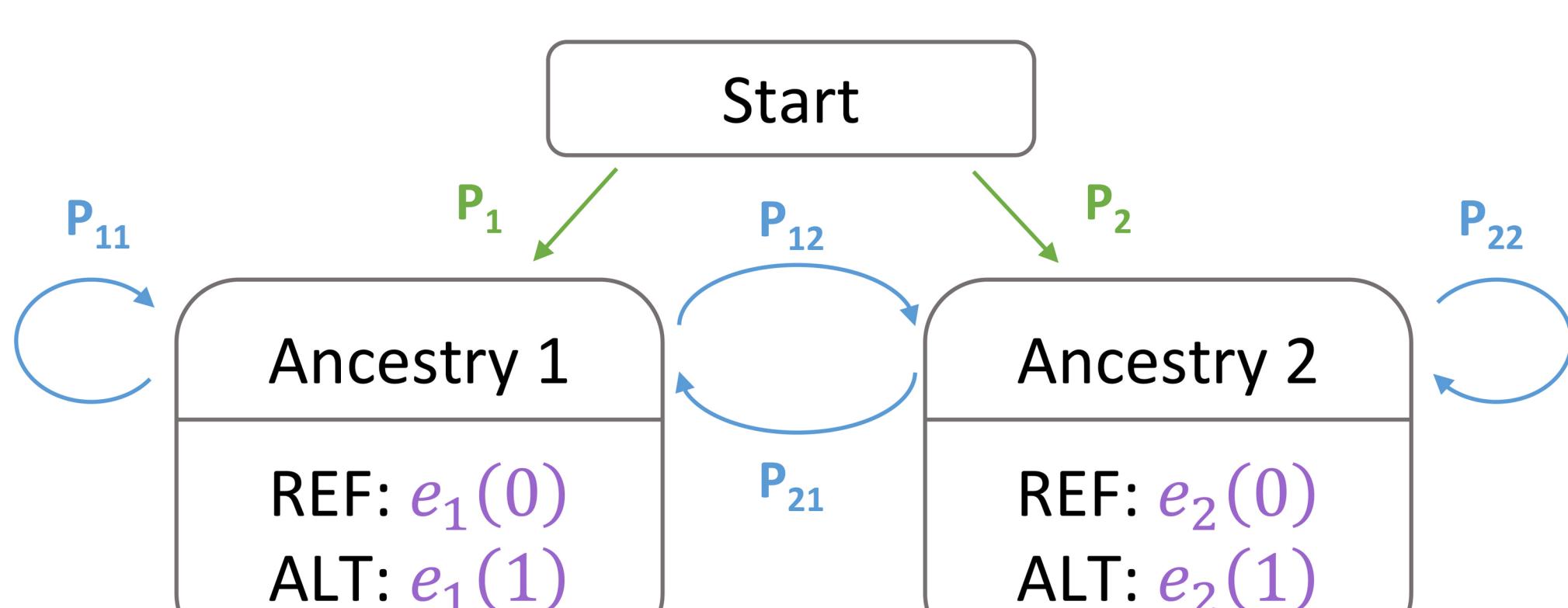
Consensus Local Ancestry Estimate

Hidden Markov Model

Example: For $l = 1$

Hidden: Sequence of local ancestry states, $S = \{s_1, s_2, \dots, s_i, \dots\}$

Observe: Sequence of haplotypes, $X = \{x_1, x_2, \dots, x_i, \dots\}$



Hidden Markov Model (cont.)

Initialize (Forward Direction):

$$P_1(x_1, 1) = \log(P_1) + \log(e_1(x_1)) \quad P_2(x_1, 1) = \log(P_2) + \log(e_2(x_1))$$

Iterate:

$$P_1(x_i, i) = \log(e_1(x_i)) + \max[\log(P_{11}) + \log(P_1(x_{i-1}, i-1)), \log(P_{21}) + \log(P_2(x_{i-1}, i-1))]$$

$$P_2(x_i, i) = \log(e_2(x_i)) + \max[\log(P_{22}) + \log(P_2(x_{i-1}, i-1)), \log(P_{12}) + \log(P_1(x_{i-1}, i-1))]$$

State selection:

If $P_1(x_i, i) > P_2(x_i, i)$, then $S_i = \text{Ancestry 1}$

If $P_2(x_i, i) > P_1(x_i, i)$, then $S_i = \text{Ancestry 2}$

Results

Figure 1: Local ancestry estimates for homologous chromosomes from 3 ASW samples. Each sample has some discrepancy between forward and backward estimates. Yet, summed local ancestry estimate close to Admixture results.

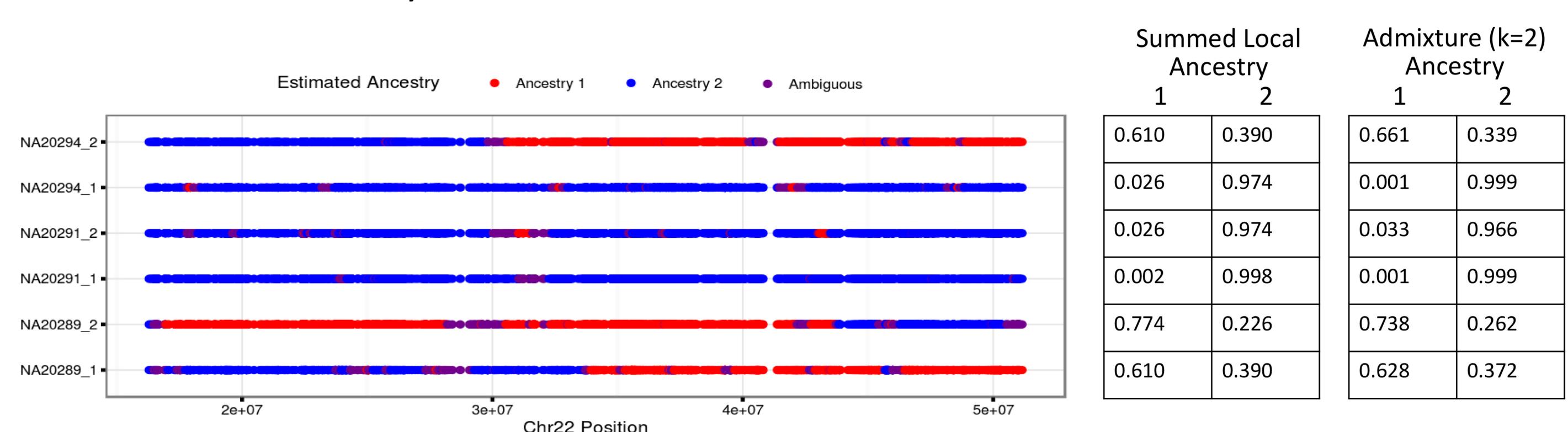


Figure 2: Summation of local ancestry estimates highly correlated with Admixture (k=2) results

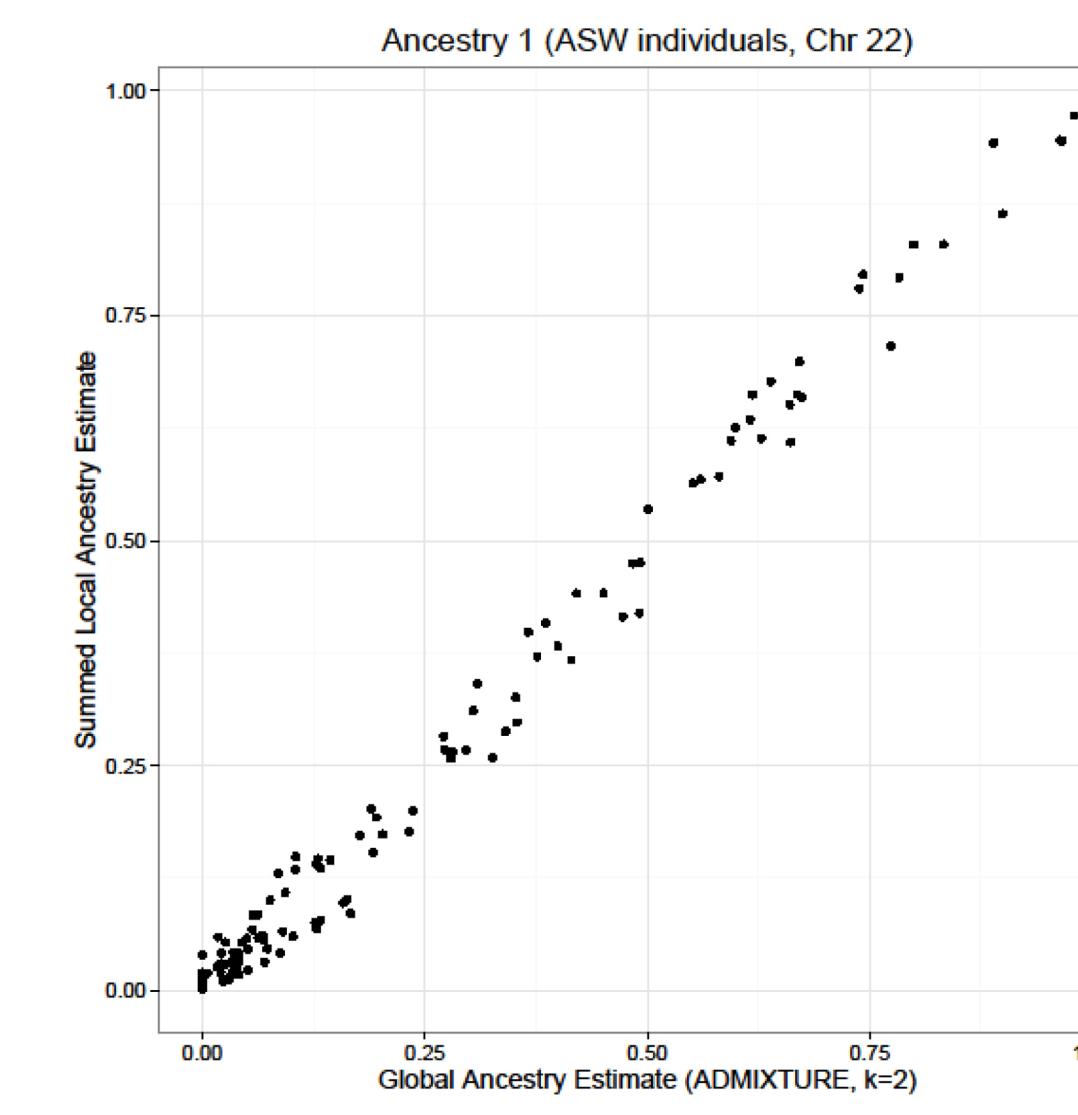
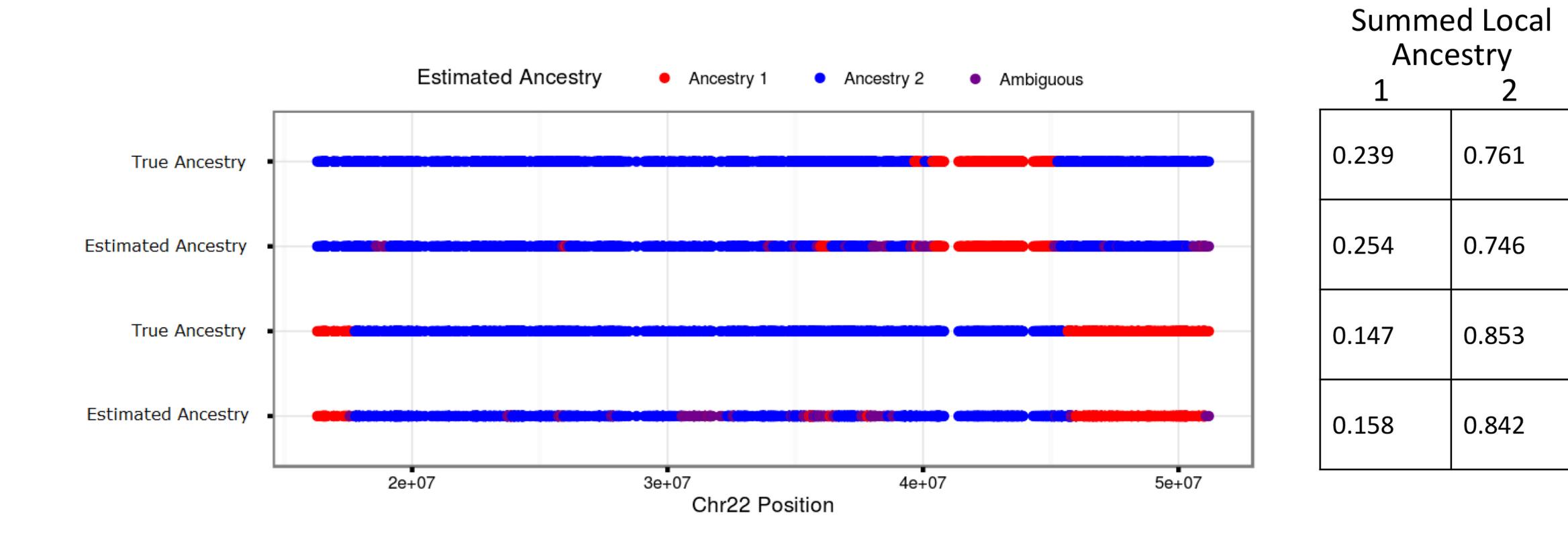


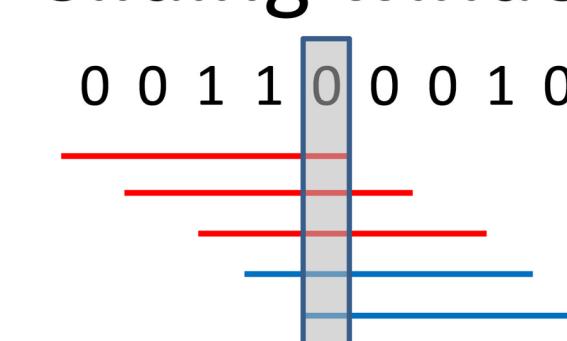
Figure 3: Local ancestry estimates for known mixture between CEU and YRI ancestry.



Future Improvements

1. Reduce short, erroneous ancestral flips

a. Sliding windows of haplotypes

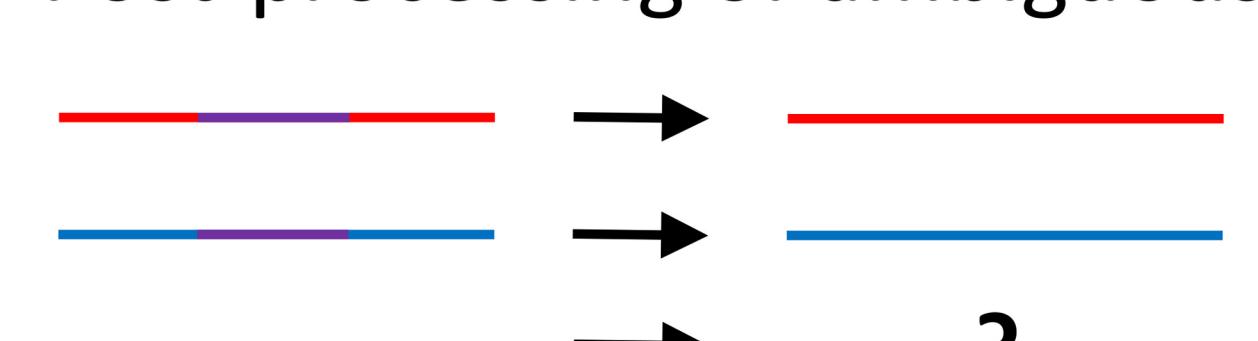


Selection based on:

- Number of windows with ancestry
- Weighted by size of haplotype
- Weighted by confidence in ancestry call

b. Incorporate recombination maps

c. Post-processing of ambiguous local ancestry estimates



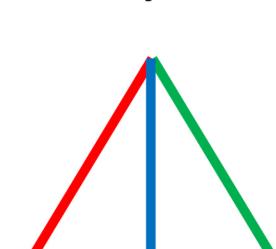
Depending on:

- Length of ambiguous region
- Relative confidence in ancestry call by direction

2. Inference with greater than 2 ancestral populations

3. Hierarchical painting of ancestry

Inference methods assume no evolutionary relationship



Evolutionary relationships do exist between populations

